

To Tune or Not to Tune the Number of Trees in Random Forest

Philipp Probst

PROBST@IBE.MED.UNI-MUENCHEN.DE

*Institut für medizinische Informationsverarbeitung, Biometrie und Epidemiologie
Marchioninstr. 15, 81377 München*

Anne-Laure Boulesteix

BOULESTEIX@IBE.MED.UNI-MUENCHEN.DE

*Institut für medizinische Informationsverarbeitung, Biometrie und Epidemiologie
Marchioninstr. 15, 81377 München*

Editor: Isabelle Guyon

Abstract

The number of trees T in the random forest (RF) algorithm for supervised learning has to be set by the user. It is unclear whether T should simply be set to the largest computationally manageable value or whether a smaller T may be sufficient or in some cases even better. While the principle underlying bagging is that more trees are better, in practice the classification error rate sometimes reaches a minimum before increasing again for increasing number of trees. The goal of this paper is four-fold: (i) providing theoretical results showing that the expected error rate may be a non-monotonous function of the number of trees and explaining under which circumstances this happens; (ii) providing theoretical results showing that such non-monotonous patterns cannot be observed for other performance measures such as the Brier score and the logarithmic loss (for classification) and the mean squared error (for regression); (iii) illustrating the extent of the problem through an application to a large number ($n = 306$) of datasets from the public database OpenML; (iv) finally arguing in favor of setting T to a computationally feasible large number as long as classical error measures based on average loss are considered.

Keywords: Random forest, number of trees, bagging, out-of-bag, error rate

1. Introduction

The random forest (RF) algorithm for classification and regression, which is based on the aggregation of a large number T of decision trees, was first described in its entirety by Breiman (2001). T is one of several important parameters which have to be carefully chosen by the user. Some of these parameters are *tuning parameters* in the sense that both too high and too low parameter values yield sub-optimal performances; see Segal (2004) for an early study on the effect of such parameters. It is unclear, however, whether the number of trees T should simply be set to the largest computationally manageable value or whether a smaller T may be sufficient or in some cases even better, in which case T should ideally be tuned carefully. This question is relevant to any user of RF and has been the topic of much informal discussion in the scientific community, but has to our knowledge never been addressed systematically from a theoretical and empirical point of view.

Breiman (2001) provides proofs of convergence for the generalization error in the case of classification random forest for growing number of trees. This means that the error rate

for a given test or training dataset converges to a certain value. Moreover, Breiman (2001) proves that there exists an upper bound for the generalization error. Similarly he proves the convergence of the mean squared generalization error for regression random forests and also provides an upper bound. However, these results do not answer the question of whether the number of trees is a tuning parameter or should be set as high as computationally feasible, although convergence properties may at first view be seen as an argument in favor of a high number of trees. Breiman (1996a) and Friedman (1997) note that bagging and aggregation methods can make good predictors better but poor predictors can be transformed into worse. Hastie et al. (2001) show in a simple example that for a single observation that is incorrectly classified (in the binary case), bagging can worsen the expected missclassification rate. In Section 3.1 we will further analyse this issue and examine the outcome of aggregating performances for several observations.

Since each tree is trained individually and without knowledge of previously trained trees, however, the risk of overfitting when adding more trees discussed by Friedman (2001) in the case of boosting is not relevant here.

The number of trees is sometimes considered as a tuning parameter in current literature (Raghu et al., 2015); see also Barman et al. (2014) for a study in which different random seeds are tested to obtain better forests—a strategy implicitly assuming that a random forest with few trees may be better than a random forest with many trees. The R package `RFmarkerDetector` (Palla and Armano, 2016) even provides a function, `'tuneNTREE'`, to tune the number of trees. Of note, the question of whether a smaller number of trees may be better has often been discussed in online forums (see Supplementary File 1 for a non-exhaustive list of links) and seems to remain a confusing issue to date, especially for beginners.

A related but different question is whether a smaller number of trees is *sufficient* (as opposed to “better”) in the sense that more trees do not improve accuracy. This question is examined, for example, in the very early study by Latinne et al. (2001) or by Hernández-Lobato et al. (2013). Another important contribution to that question is the study by Oshiro et al. (2012), which compared the performance in terms of the Area Under the ROC Curve (AUC) of random forests with different numbers of trees on 29 datasets. Their main conclusion is that the performance of the forest does not always substantially improve as the number of trees grows and after having trained a certain number of trees (in their case 128) the AUC performance gain obtained by adding more trees is minimal. The study of Oshiro et al. (2012) provides important empirical support for the existence of a “plateau”, but does not directly address the question of whether a smaller number of trees may be substantially better and does not investigate this issue from a theoretical perspective, thus making the conclusions dependent on the 29 examined datasets.

In this context, the goal of our paper is four-fold: (i) providing theoretical results showing that, in the case of binary classification, the expected error rate may be a non-monotonous function of the number of trees and explaining under which circumstances this happens; (ii) providing theoretical results showing that such non-monotonous patterns cannot be observed for other performance measures such as the Brier score and the logarithmic loss (for classification) and the mean squared error (for regression); (iii) illustrating the extent of the problem through an application to a large number ($n = 306$) of datasets from the public database OpenML; (iv) finally arguing in favor of setting it to a computationally feasible

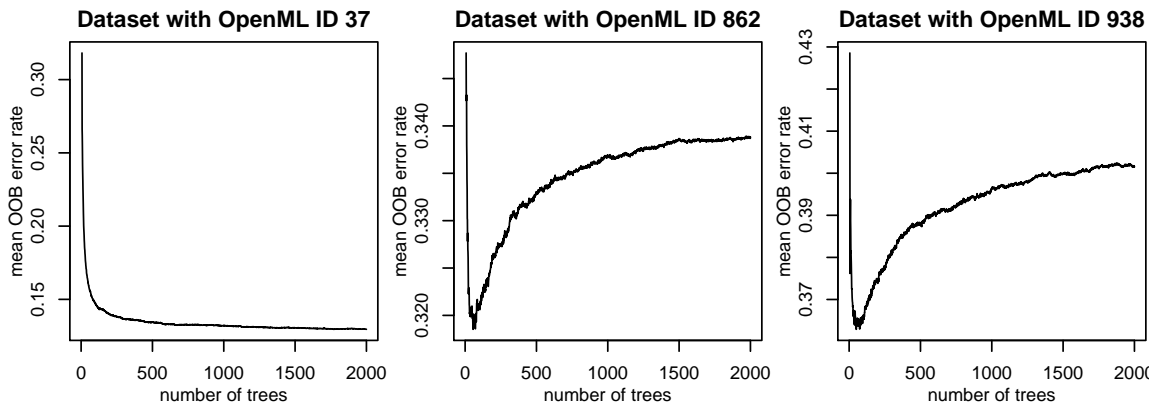


Figure 1: Mean OOB error rate curves for OpenML datasets with IDs 37, 862 and 938. The curves are averaged over 1000 independent runs of random forest.

large number as long as classical error measures based on average loss are considered. Furthermore, we introduce our new R package `OOBCurve`, which can be used to examine the convergence of various performance measures.

To set the scene, we first address this issue empirically by looking at the curve depicting the out-of-bag (OOB) error rate (see Section 2 for a definition of the OOB error) for different number of trees (also called OOB error rate curve) for various datasets from the OpenML database (Vanschoren et al., 2013). To obtain more stable results and better estimations for the expected error rate we repeat this procedure 1000 times for each dataset and average the results.

For most datasets we observe monotonously decreasing curves with growing number of trees as in the left panel of Figure 1, while others yield strange non-monotonous patterns, for example the curves of the datasets with the OpenML ID 862 and 938, which are also depicted in Figure 1. The initial error rate drops steeply before starting to increase after a certain number of trees before finally reaching a plateau.

At first view, such non-monotonous patterns are a clear argument in favor of tuning T . We claim, however, that it is important to understand why and in which circumstances such patterns happen in order to decide whether or not T should be tuned in general. In Section 3, we address this issue from a theoretical point of view, by formulating the expected error rate as a function of the probabilities ε_i of correct classification by a single tree for each observation i of the training dataset, for $i = 1, \dots, n$ (with n denoting the size of the training dataset). This theoretical view provides a clear explanation of the non-monotonous error rate curve patterns in the case of classification. With a similar approach, we show that such non-monotonous patterns cannot be obtained with the Brier score or the logarithmic loss as performance measures, which are based on probability estimations and also not for the mean squared error in the case of regression. Only for the AUC we can see non-monotonous curves as well.

The rest of this paper is structured as follows. Section 2 gives a brief introduction into random forest and performance estimation. Theoretical results are presented in Section 3, while the results of a large empirical study based on 306 datasets from the public database OpenML are reported in Section 4. More precisely, we empirically validate our theoretical model for the error as a function of the number of trees as well as our statements regarding the properties of datasets yielding non-monotonous patterns. We finally argue in Section 5 that there is no inconvenience—except additional computational cost—in adding trees to a random forest and that T should thus not be seen as a tuning parameter as long as classical performance measures based on the average loss are considered.

2. Background: Random Forest and Measures of Performance

In this section we introduce the random forest method, the general notation and some well known performance measures.

2.1 Random Forest

The random forest (RF) is an ensemble learning technique consisting of the aggregation of a large number T of decision trees, resulting in a reduction of variance compared to the single decision trees. In this paper we consider the original version of RF first described by Breiman (2001), while acknowledging that other variants exist, for example RF based on conditional inference trees (Hothorn et al., 2006) which address the problem of variable selection bias investigated by Strobl et al. (2007). Our considerations are however generalizable to many of the available RF variants and other methods that use randomization techniques.

A prediction is obtained for a new observation by aggregating the predictions made by the T single trees. In the case of regression RF, the most straightforward and common procedure consists of averaging the prediction of the single trees, while majority voting is usually applied to aggregate classification trees. This means that the new observation is assigned to the class that was most often predicted by the T trees.

While RF can be used for various types of response variables including censored survival times or (as empirically investigated in Section 4) multicategorical variables, in this paper we mainly focus on the two most common cases, binary classification and regression.

2.2 General Notations

From now on, we consider a fixed training dataset D consisting of n observations, which is used to derive prediction rules by applying the RF algorithm with a number T of trees. Ideally, the performance of these prediction rules is estimated based on an independent test dataset, denoted as D_{test} , consisting of n_{test} test observations.

Considering the i th observation from the test dataset ($i = 1, \dots, n_{test}$), we denote its true response as y_i , which can be either a numeric value (in the case of regression) or the binary label 0 vs. 1 (in the case of binary classification). The predicted value output by tree t (with $t = 1, \dots, T$) is denoted as \hat{y}_{it} , while \hat{y}_i stands for the predicted value output by the whole random forest. Note that, in the case of regression, \hat{y}_i is usually obtained by

averaging as

$$\hat{y}_i = \frac{1}{T} \sum_{t=1}^T \hat{y}_{it}.$$

In the case of classification, \hat{y}_i is usually obtained by majority voting. For binary classification, it is equivalent to computing the same average as for regression, which now takes the form

$$\hat{p}_i = \frac{1}{T} \sum_{t=1}^T I(\hat{y}_{it} = 1)$$

and is denoted as \hat{p}_i (standing for probability), and finally deriving \hat{y}_i as

$$\hat{y}_i = \begin{cases} 1 & \text{if } \hat{p}_i > 0.5, \\ 0 & \text{otherwise.} \end{cases}$$

2.3 Measures of Performance for Binary Classification and Regression

In regression as well as in classification, the performance of a RF for observation i is usually quantified through a so-called loss function measuring the discrepancy between the true response y_i and the predicted response \hat{y}_i or, in the case of binary classification, between y_i and \hat{p}_i . For both regression and binary classification, the classical and most straightforward measure is defined for observation i as

$$e_i = (y_i - \hat{y}_i)^2 = L(y_i, \hat{y}_i),$$

with $L(.,.)$ standing for the loss function $L(x, y) = (x - y)^2$. In the case of regression this is simply the squared error. Another common loss function in the regression case is the absolute loss $L(x, y) = |x - y|$. For binary classification both measures simplify to $e_i = 0$ if observation i is classified correctly by the RF, $e_i = 1$ otherwise, which we will simply denote as *error* from now on. One can also consider the performance of single trees, that means the discrepancy between y_i and \hat{y}_{it} . We define e_{it} as

$$e_{it} = L(y_i, \hat{y}_{it}) = (y_i - \hat{y}_{it})^2$$

and the mean error—a quantity we need to derive our theoretical results on the dependence of performance measures on the number of tree T —as

$$\varepsilon_i = E(e_{it}),$$

where the expectation is taken over the possible trees conditionally on D . The term ε_i can be interpreted as the difficulty to predict y_i with single trees. In the case of binary classification, we have $(y_i - \hat{y}_{it})^2 = |y_i - \hat{y}_{it}|$ and ε_i can be simply estimated as $|y_i - \hat{p}_i|$ from a RF with a large number of trees.

In the case of binary classification, it is also common to quantify performance through the use of the Brier score, which has the form

$$b_i = (y_i - \hat{p}_i)^2 = L(y_i, \hat{p}_i)$$

or of the logarithmic loss

$$l_i = -(y_i \ln(\hat{p}_i) + (1 - y_i) \ln(1 - \hat{p}_i)).$$

Both of them are based on \hat{p}_i rather than \hat{y}_i , and can thus be only defined for the whole RF and not for single trees.

The area under the ROC curve (AUC) cannot be expressed in terms of single observations, as it takes into account all observations at once by ranking the \hat{p}_i -values. It can be interpreted as the probability that the classifier ranks a randomly chosen observation with $y_i = 1$ higher than a randomly chosen observation with $y_i = 0$. The larger the AUC, the better the discrimination between the two classes. The (empirical) AUC is defined as

$$\text{AUC} = \frac{\sum_{i=1}^{n_1} \sum_{j=1}^{n_2} S(\hat{p}_i^*, \hat{p}_j^{**})}{n_1 n_2},$$

where $\hat{p}_1^*, \dots, \hat{p}_{n_1}^*$ are probability estimations for the n_1 observations with $y_i = 1$, $\hat{p}_1^{**}, \dots, \hat{p}_{n_2}^{**}$ are probability estimations for the n_2 observations with $y_i = 0$ and $S(., .)$ is defined as $S(p, q) = 0$ if $p < q$, $S(p, q) = 0.5$ if $p = q$ and $S(p, q) = 1$ if $p > q$. The AUC can also be interpreted as the Mann-Whitney U-Statistic divided by the product of n_1 and n_2 .

2.4 Measures for Multiclass Classification

The measures defined in the previous section can be extended to the multiclass classification case. Let K denote the number of classes ($K > 2$). The response y_i takes values in $\{1, \dots, K\}$. The error for observation i is then defined as

$$e_i = I(y_i \neq \hat{y}_i).$$

We denote the estimated probability of class k for observation i as

$$\hat{p}_{ik} = \frac{1}{T} \sum_{t=1}^T I(\hat{y}_{it} = k).$$

The logarithmic loss is then defined as

$$l_i = \sum_{k=1}^K -I(y_i = k) \log(\hat{p}_{ik})$$

and the generalized Brier score is defined as

$$b_i = \sum_{k=1}^K (\hat{p}_{ik} - I(y_i = k))^2,$$

which in the binary case is twice the value of the definition that was used in the previous section. Following Hand and Till (2001), the AUC can also be generalized to the multiclass case as

$$\text{AUC} = \frac{1}{K(K-1)} \sum_{j=1}^K \sum_{\substack{k=1 \\ k \neq j}}^K \text{AUC}(j, k),$$

where $\text{AUC}(j, k)$ is the AUC between class k and j , see also Ferri et al. (2009) for more details. It is equivalent to the definition given in Section 2.3 in the binary classification case.

2.5 Test Dataset Error vs. Out-of-Bag Error

In the cases where a test dataset D_{test} is available, performance can be assessed by averaging the chosen performance measure (as described in the previous paragraphs) over the n_{test} observations. For example the classical error rate (for binary classification) and the mean squared error (for regression) are computed as

$$\frac{1}{n_{test}} \sum_{i=1}^{n_{test}} L(y_i, \hat{y}_i),$$

with $L(x, y) = (x - y)^2$, while the mean absolute error (for regression) is obtained by defining $L(., .)$ as $L(x, y) = |x - y|$. Note that, in the context of regression, Rousseeuw (1984) proposes to consider the median $med(L(y_1, \hat{y}_1), \dots, L(y_{n_{test}}, \hat{y}_{n_{test}}))$, instead of averaging, which results in the median squared error for the loss function $L(x, y) = (x - y)^2$ and in the median absolute error for the loss function $L(x, y) = |x - y|$. These measures are more robust against outliers and contamination (Rousseeuw, 1984).

An alternative to the use of a test dataset is the out-of-bag error which is calculated by using the out-of-bag (OOB) estimations of the training observations. OOB predictions are calculated by predicting the class, the probability (in the classification case) or the real value (in the regression case) for each training observation i (for $i = 1, \dots, n$) by using only the trees for which this observation was not included in the bootstrap sample (i.e., it was not used to construct the tree). Note that these predictions are obtained based on a subset of trees—including on average $T \times 0.368$ trees. These predictions are ultimately compared to the true values by calculating performance measures (see Sections 2.3, 2.4 and 2.5).

3. Theoretical Results

In this section we compute the expected performance—according to the error, the Brier score and the logarithmic loss outlined in Section 2.3—of a binary classification or regression RF consisting of T trees as estimated based on the n_{test} test observations, while considering the training dataset as fixed. For the AUC we prove that it can be a non-monotonous function in T . The case of other measures (mean absolute error, median of squared error and median of absolute error for regression) and multiclass classification is much more more complex to investigate from a theoretical point of view. It will be examined empirically in Section 4.

In this section we are concerned with *expected* performances, where expectation is taken over the sets of T trees. Our goal is to study the monotonicity of the expected errors with respect to T . The number T of trees is considered a parameter of the RF and now mentioned in parentheses everytime we refer to the whole forest.

3.1 Error Rate (Binary Classification)

We first show that for single observations the expected error rate curve can be increasing and then show exemplified how this can influence the shape of the average curve of several observations. The observation that bagging can worsen the expected error rate of a single observation was already done by Hastie et al. (2001), Breiman (1996a) and Friedman (1997). In this section we provide a general formula explaining this observation, and then extend our theoretical considerations to further performance measures in the following sections.

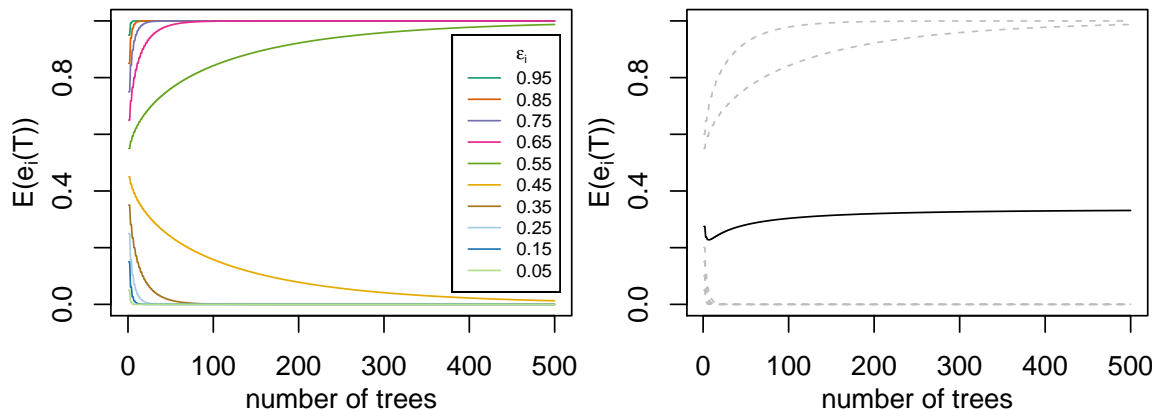


Figure 2: Left: Expected error rate curves for different ε_i values. Right: Plot of the average curve (black) of the curves with $\varepsilon_1 = 0.05$, $\varepsilon_2 = 0.1$, $\varepsilon_3 = 0.15$, $\varepsilon_4 = 0.2$, $\varepsilon_5 = 0.55$ and $\varepsilon_6 = 0.6$ (depicted in grey and dotted)

3.1.1 THEORETICAL CONSIDERATIONS

Let us first consider the classical error rate $e_i(T)$ for observation i with a RF including T trees and derive its expectation, conditionally on the training set D ,

$$E(e_i(T)) = E\left(I\left(\frac{1}{T} \sum_{t=1}^T e_{it} > 0.5\right)\right) = P\left(\sum_{t=1}^T e_{it} > 0.5 \cdot T\right).$$

We note that e_{it} is a binary variable with $E(e_{it}) = \varepsilon_i$. Given a fixed training dataset D and observation i , the e_{it} , $t = 1, \dots, T$ are mutually independent. It follows that the sum $X_i = \sum_{t=1}^T e_{it}$ follows the binomial distribution $B(T, \varepsilon_i)$. It is immediate that the contribution of observation i to the expected error rate, $P(X_i > 0.5 \cdot T)$, is an increasing function in T for $\varepsilon_i > 0.5$ and a decreasing function in T for $\varepsilon_i < 0.5$.

Note that so far we ignored the case where $\sum_{t=1}^T e_{it} = 0.5 \cdot T$, which may happen when T is even. In this case, the standard implementation in R (`randomForest`) assigns the observation randomly to one of the two classes. This implies that $0.5 \cdot P(\sum_{t=1}^T e_{it} = 0.5 \cdot T)$ has to be added to the above term, which does not affect our considerations on the ε_i 's role.

3.1.2 IMPACT ON ERROR RATE CURVES

The error rate curve for observation i is defined as the curve described by the function $e_i : T \rightarrow \mathbb{R}$. The expectation $E(e_i(T))$ of the error rate curve for observation i with the mentioned adjustment in the case of an even number of trees can be seen in the left plot of Figure 2 for different values of ε_i . Very high and very low values of ε_i lead to rapid convergence, while for ε_i -values close to 0.5 more trees are needed to reach the plateau. The error rate curve obtained for a test dataset consists of the average of the error rate curves of the single observations. Of course, if trees are good classifiers we should have $\varepsilon_i < 0.5$ for most observations. In many cases, observations with $\varepsilon_i > 0.5$ will be compensated by

observations with $\varepsilon_i < 0.5$ in such a way that the expected error rate curve is monotonously decreasing. This is typically the case if there are many observations with $\varepsilon_i \approx 0$ and a few with $\varepsilon_i \approx 1$. However, if there are many observations with $\varepsilon_i \approx 0$ and a few observations with $\varepsilon_i \geq 0.5$ that are close to 0.5, the expected error rate curve initially falls down quickly because of the observation with $\varepsilon_i \approx 0$ and then grows again slowly as the number of trees increases because of the observations with $\varepsilon_i \geq 0.5$ close to 0.5. In the right plot of Figure 2 we can see (black solid line) the mean of the expected error rate curves for $\varepsilon_1 = 0.05$, $\varepsilon_2 = 0.1$, $\varepsilon_3 = 0.15$, $\varepsilon_4 = 0.2$, $\varepsilon_5 = 0.55$ and $\varepsilon_6 = 0.6$ (displayed as gray dashed lines) and can see exactly the non-monotonous pattern that we expected: due to the ε_i 's 0.55 and 0.6 the average curve increases again after reaching a minimum. In Section 4 we will see that the two example datasets whose non-monotonous out-of-bag error rate curves are depicted in the introduction have a similar distribution of ε_i .

We see that the convergence rate of the error rate curve is only dependent on the distribution of the ε_i 's of the observations. Hence, the convergence rate of the error rate curve is not directly dependent on the number of observations n or the number of features, but these characteristics could influence the empirical distribution of the ε_i 's and hence possibly the convergence rate as outlined in Section 4.4.1.

3.2 Brier Score (Binary Classification) and Squared Error (Regression)

We now turn to the Brier score and compute the expected Brier score contribution of observation i for a RF including T trees, conditional on the training set D . We obtain

$$\begin{aligned} E(b_i(T)) &= E((y_i - \hat{p}_i(T))^2) = E\left(\left(y_i - \frac{1}{T} \sum_{t=1}^T \hat{y}_{it}\right)^2\right) \\ &= E\left(\left(\frac{1}{T} \sum_{t=1}^T (y_i - \hat{y}_{it})\right)^2\right) = E\left(\left(\frac{1}{T} \sum_{t=1}^T e_{it}\right)^2\right). \end{aligned}$$

From $E(Z^2) = E(Z)^2 + Var(Z)$ with $Z = \frac{1}{T} \sum_{t=1}^T e_{it}$ it follows:

$$E(b_i(T)) = E(e_{it})^2 + \frac{Var(e_{it})}{T},$$

which is obviously a strictly monotonous decreasing function of T . This also holds for the average over the observations of the test dataset. In the case of binary classification, we have $e_{it} \sim \mathcal{B}(1, \varepsilon_i)$, yielding $E(e_{it}) = \varepsilon_i$ and $Var(e_{it}) = \varepsilon_i(1 - \varepsilon_i)$, thus allowing the formulation of $E(b_i(T))$ as $E(b_i(T)) = \varepsilon_i^2 + \frac{\varepsilon_i(1 - \varepsilon_i)}{T}$. Note that the formula $E(b_i(T)) = E(e_{it})^2 + Var(e_{it})/T$ is also valid for the squared error in the regression case, except that in this case we would write \hat{y}_i instead of \hat{p}_i in the first line.

3.3 Logarithmic Loss (Binary Classification)

As outlined in Section 2.3, another usual performance measure based on the discrepancy between y_i and \hat{p}_i is the logarithmic loss $l_i(T) = -(y_i \ln(\hat{p}_i(T)) + (1 - y_i) \ln(1 - \hat{p}_i(T)))$. Noticing that $\hat{p}_i(T) = 1 - \frac{1}{T} \sum_{t=1}^T e_{it}$ for $y_i = 1$ and $\hat{p}_i(T) = \frac{1}{T} \sum_{t=1}^T e_{it}$ for $y_i = 0$, it can

be in both cases $y_i = 0$ and $y_i = 1$ reformulated as

$$l_i(T) = -\ln \left(1 - \frac{1}{T} \sum_{t=1}^T e_{it} \right).$$

In the following we ensure that the term inside the logarithm is never zero by adding a very small value a to $1 - \frac{1}{T} \sum_{t=1}^T e_{it}$. The logarithmic loss $l_i(T)$ is then always defined and its expectation exists. This is similar to the solution adopted in the `mlr` package, where 10^{-15} is added in case that the inner term of the logarithm equals zero.

With $Z := 1 - \frac{1}{T} \sum_{t=1}^T e_{it} + a$, we can use the Taylor expansion,

$$\begin{aligned} E[f(Z)] &= E[f(\mu_Z + (Z - \mu_Z))] \\ &\approx E \left[f(\mu_Z) + f'(\mu_Z)(Z - \mu_Z) + \frac{1}{2} f''(\mu_Z)(Z - \mu_Z)^2 \right] \\ &= f(\mu_Z) + \frac{f''(\mu_Z)}{2} \cdot \text{Var}(Z) = f(E(Z)) + \frac{f''(E(Z))}{2} \cdot \text{Var}(Z) \end{aligned}$$

where μ_Z stands for $E(Z)$ and $f(\cdot)$ as $f(\cdot) = -\ln(\cdot)$. We have $\text{Var}(Z) = \frac{\varepsilon_i(1-\varepsilon_i)}{T}$, $E(Z) = 1 - \varepsilon_i + a$, $f(E(Z)) = -\ln(1 - \varepsilon_i + a)$ and $f''(E(Z)) = (1 - \varepsilon_i + a)^{-2}$, finally yielding

$$E(l_i(T)) \approx -\ln(1 - \varepsilon_i + a) + \frac{\varepsilon_i(1 - \varepsilon_i)}{2T(1 - \varepsilon_i + a)^2},$$

which is obviously a decreasing function of T . The Taylor approximation gets better and better for increasing T , since the variance of $l_i(T)$ decreases with increasing T and thus $l_i(T)$ tends to get closer to its expectancy.

3.4 Area Under the ROC Curve (AUC) (Classification)

For the AUC, considerations such as those we made for the error rate, the Brier score and the logarithmic loss are impossible, since the AUC is not the sum of individual contributions of the observations. It is however relatively easy to see that the expected AUC is not always an increasing function of the number T of trees. For example, think of the trivial example of a test dataset consisting of two observations with responses y_1 resp. y_2 and $E(\hat{p}_1(T)) = 0.4$ resp. $E(\hat{p}_2(T)) = 0.6$. If $y_1 = 0$ and $y_2 = 1$, the expected AUC curve increases monotonously with T , as the probability of a correct ordering according to the calculated scores $\hat{p}_1(T)$ and $\hat{p}_2(T)$ increases. However, if $y_1 = 1$ and $y_2 = 0$, we obtain a monotonously decreasing function, as the probability of a wrong ordering gets higher with increasing number of trees. It is easy to imagine that for different combinations of $E(\hat{p}_i(T))$, one can obtain increasing curves, decreasing curves or non-monotonous curves.

3.5 Adapting the Models to the OOB Error

The ‘‘OOB estimator’’ of the performance outlined in Section 2.5 is commonly considered as an acceptable proxy of the performance estimator obtained through the use of an independent test dataset or through resampling-techniques such as cross-validation (Breiman, 1996b) for a random forest including $T \times 0.368$ trees. Compared to these techniques, the

OOB estimator has the major advantage that it neither necessitates to fit additional random forests (which is advantageous in terms of computational resources) nor to reduce the size of the dataset through data splitting. For these reasons, we will consider OOB performance estimators in our empirical study.

However, if we consider the OOB error instead of the test error from an independent dataset, the formulas given in the previous subsections are not directly applicable. After having trained T trees, for making an OOB estimation for an observation we can only use the trees for which the observation was out-of-bag. If we take a simple bootstrap sample from the n training observation when bagging we have *on average* only $T \cdot (1 - \frac{1}{n})^n \approx T \cdot \exp(-1) \approx T \cdot 0.368$ trees for predicting the considered observation. This means that we would have to replace T by $T \cdot \exp(-1)$ in the above formulas and that the formulas are no longer exact because $T \cdot \exp(-1)$ is only an average. Nonetheless it is still a good approximation as confirmed in our benchmark experiments.

4. Empirical Results

This section shows a large-scale empirical study based on 193 classification tasks and 113 regression tasks from the public database OpenML (Vanschoren et al., 2013). The datasets are downloaded with the help of the `OpenML R` package (Casalicchio et al., 2017). The goals of this study are to (i) give an order of magnitude of the frequency of non-monotonous patterns of the error rate curve in real data settings; (ii) empirically confirm our statement that observations with ε_i greater than (but close to) 0.5 are responsible for non-monotonous patterns; (iii) analyse the results for other classification measures, the multiclass classification and several regression measures; (iv) analyse the convergence rate of the OOB curves.

4.1 Selection of Datasets

To select the datasets to be included in our study we define a set of candidate datasets—in our case the datasets available from the OpenML platform (Vanschoren et al., 2013)—and a set of inclusion criteria as recommended in Boulesteix et al. (2017). In particular, we do not select datasets with respect to the results they yield, thus warranting representativity.

Our inclusion criteria are as follows: (i) the dataset has predefined tasks in OpenML (see Vanschoren et al., 2013, for details on the OpenML nomenclature); (ii) it includes less than 1000 observations; (iii) it includes less than 1000 features. The two latter criteria aim at keeping the computation time feasible.

Cleaning procedures such as the deletion of duplicated datasets (whole datasets that appear twice in the OpenML database) are also applied to obtain a decent collection of datasets. No further modification of the tasks and datasets were done.

This procedure yields a total of 193 classification tasks and 113 regression tasks.

From the 193 classification tasks, 149 are binary classification tasks and 44 multiclass classification tasks.

The tasks contained easy, medium and difficult tasks - for binary classification tasks the mean (out-of-bag) AUC of a random forest with 2000 trees was 0.841, the minimum 0.502, the first quartile 0.732, the median 0.870, the third quartile 0.962 and the maximum 1. Similarly the regression tasks contained easy and difficult tasks with a mean R^2 of 0.559.

4.2 Study Design

For each dataset we run the RF algorithm with $T = 2000$ trees 1000 times successively with different seeds using the R package `randomForest` (Liaw and Wiener, 2002) with the default parameter settings. We choose 2000 trees because in a preliminary study on a subset of the datasets we could observe that convergence of the OOB curves was reached within these 2000 trees. Note that all reported results regarding the performance gain and convergence are made with the out-of-bag predictions. As for these predictions on average only $\exp(-1) \cdot T$ of the T trees are used, the convergence of independent test data is faster by the factor 2.7. For the classification tasks we calculate the OOB curves for the error rate, the balanced error rate, the (multiclass) Brier score, the logarithmic loss and the (multiclass) AUC using our new package `OOBCurve`, see details in the next section.

For the regression tasks we calculate the OOB curves using the mean squared error, the mean absolute error, the median squared error and the median of absolute error as performance measures. We parallelize the computations using the R package `batchtools` (version 0.9.0) (Lang et al., 2017). For each measure and each dataset, the final curve is obtained by simply averaging over the 1000 runs of RF. We plot each of them in three files separately for binary classification, multiclass classification and regression. In the plots the x-axis starts at $T = 11$ since overall performance estimates are only defined if each observation was out-of-bag in at least one of the T trees, which is not always the case in practice for $T < 10$. We plot the curves only until $T = 500$, as no interesting patterns can be observed after this number of trees (data not shown). The graphics, the R-codes and the results of our experiment can be found on <https://github.com/PhilippPro/tuneNtree>.

4.3 The R Package `OOBCurve`

The calculation of out-of-bag estimates for different performance measures is implemented in our new R package `OOBCurve`. More precisely, it takes a random forest constructed with the R package `randomForest` (Liaw and Wiener, 2002) or `ranger` (Wright, 2016) as input and can calculate the OOB curve for any measure that is available from the `mlr` package (Bischl et al., 2016). The `OOBCurve` package is available on CRAN R package repository and also on Github (<https://github.com/PhilippPro/OOBCurve>). It is also possible to calculate OOB curves of other hyperparameters of RF such as `mtry` with this package.

4.4 Results for Binary Classification

The average gain in performance in the out-of-bag performance for 2000 trees instead of 11 trees is -0.0324 for the error rate, -0.0683 for the brier score, -2.383 for the logarithmic loss and 0.0553 for the AUC. In the following we will concentrate on the visual analysis of the graphs and are especially interested in the results of the error rate.

4.4.1 OVERALL RESULTS FOR THE OOB ERROR RATE CURVES

We observe in the graphs of the OOB error rate curves that for most datasets the curve is quickly decreasing until it converges to a dataset-specific plateau value. In 16 cases which make approximately 10% of the datasets, however, the curve grows again after reaching its lowest value, leading to a value at 2000 trees that is by at least 0.005 bigger than the

lowest value of the OOB error rate curve for $T \in [10, 250]$. This happens mainly for smaller datasets, where a few observations can have a high impact on the error curve. Of these 16 cases 15 belong to the smaller half of the datasets—ordered by the number of observations multiplied with the number of features. The mean increase of these 16 datasets was 0.020 (median: 0.012). The difference in mean and median is mainly caused by one outlier where the increase was around 0.117.

4.4.2 DATASETS WITH NON-MONOTONOUS OOB ERROR RATE CURVE

We now examine in more detail the datasets yielding non-monotonous patterns. In particular, the histograms of the estimates $\hat{\varepsilon}_i = |y_i - \hat{p}_i|$ of the observation-specific errors ε_i are of interest, since our theoretical results prove that the distribution of the ε_i determines the form of the expected error rate curve. To get these histograms we compute the estimates $\hat{\varepsilon}_i$ of the observation-specific errors ε_i (as defined in Section 2.3) from a RF with a big number $T = 100000$: the more trees, the more accurate the estimates of ε_i .

The histograms for the exemplary datasets considered in the introduction (see Figure 1) are displayed in Figure 3. A typical histogram for an OOB curve with monotonously decreasing error rate curve is displayed in the left panel. The heights of the bins of this histogram of the $\hat{\varepsilon}_i$ are monotonously decreasing from 0 to 1.

The histograms for the non-monotonous error rate curves from the introduction can be seen in the middle (OpenML ID 862) and right (OpenML ID 938) panels of Figure 3. In both cases we see that a non-negligible proportion of observations have ε_i larger than but close to 0.5. This is in agreement with our theoretical results. With growing number of trees the chance that these observations are incorrectly classified increases, while the chance for observations with $\varepsilon_i \approx 0$ is already very low—and thus almost constant. Intuitively we expect such shapes of histograms for datasets with few observations—where by chance the shape of the histogram of the $\hat{\varepsilon}_i$ could look like in our two examples. For bigger datasets we expect smoother shapes of the histogram, yielding strictly decreasing error rate curves.

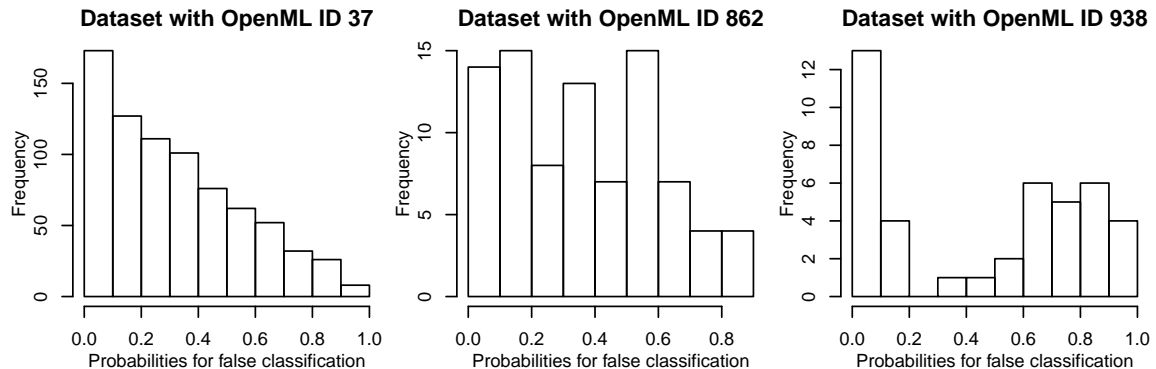


Figure 3: Histograms of the estimates of ε_i ($i = 1, \dots, n$) from random forests with 100000 trees for dataset with IDs 36, 862 and 938

	error rate		Brier score		logarithmic loss		AUC	
error rate	1.00	(1.00)	0.28	(0.44)	0.27	(0.45)	-0.18	(-0.43)
Brier score	0.72	(0.86)	1.00	(1.00)	0.96	(0.98)	-0.63	(-0.87)
logarithmic loss	0.65	(0.84)	0.93	(0.95)	1.00	(1.00)	-0.63	(-0.87)
AUC	-0.64	(-0.85)	-0.84	(-0.95)	-0.81	(-0.92)	1.00	(1.00)

Table 1: Linear (bottom-left) and rank (top-right) correlation results for binary classification datasets and for multiclass classification (in brackets)

4.4.3 OTHER MEASURES

For the Brier score and the logarithmic loss we observe, as expected, monotonically decreasing curves for all datasets. The expected AUC curve usually appears as a growing function in T . In a few datasets such as the third binary classification example (OpenML ID 905), however, it falls after reaching a maximum.

To assess the similarity between the different curves, we calculate the Bravais-Pearson linear correlation and Kendall’s τ rank correlation between the values of the OOB curves of the different performance measures and average these correlation matrices over all datasets. Note that we do not perform any correlation tests, since the assumption of independent identically distributed observations required by these tests is not fulfilled: our correlation analyses are meant to be explorative. The results can be seen in Table 1. The Brier score and logarithmic loss have the highest correlation. They are also more correlated to the AUC than to the error rate, which has the lowest correlation to all other measures.

4.5 Results for Multiclass Classification

The average gain in out-of-bag performance for 2000 trees instead of 11 trees is -0.0753 for the error rate, -0.1282 for the brier score, -5.3486 for the logarithmic loss and 0.0723 for the AUC. These values are higher than the ones from binary classification. However, the visual observations we made for the binary classification also hold for the multiclass classification. For 5 of the 44 datasets the minimum error rate for $T \in [11; 250]$ is lower by more than 0.005 than the error rate for $T = 2000$. In contrast to the binary classification case, 3 of these 5 datasets belong to the bigger half of the datasets. The results for the correlation are quite similar, although the correlation (see Table 1) is in general slightly higher than in the binary case.

4.6 Results for Regression

The average performance gain regarding the out-of-bag performance of the R^2 for 2000 trees compared to 11 trees is 0.1249. In the OOB curves for regression we can observe the monotonously decreasing pattern expected from theory in the case of the most widely used mean squared error (mse). The mean absolute error (mae) is also strictly decreasing for all the datasets considered in our study.

For the median squared error (medse) and the median absolute error (medae), we get a performance gain by using 2000 trees instead of 10 in most but not all cases (around 80% of the datasets). In many cases (around 50%) the minimum value for $T \in [11; 250]$ is smaller

than the value for $T = 2000$ which means that growing more trees is rather disadvantageous in these cases in terms of medse and medae. This could be explained by the fact that each tree in a random forest tries to minimize the squared error in the splits and therefore adding more trees to the forest will improve the mean squared error but not necessarily measures that use the median. More specifically, one could imagine that the additional trees focus on the reduction of the error for outlying observations at the price of an increase of the median error. In a simulated dataset (linear model with 200 observations, 5 relevant features and 5 non-relevant features drawn from a multivariate normal distribution) we could observe this pattern (data not shown). Without outlier all expected curves are strictly decreasing. When adding an outlier (changing the outcome of one observation to a very big value) the expected curves of mse and mae are still strictly decreasing, while the expected curves of medse and medae show are increasing for higher T . The curves of the measures which take the mean of the losses of all observations have a high linear and rank correlation (> 0.88), as well as the curves of the measures which take the median of the losses (> 0.97). Correlation between these two groups of measures are lower, around 0.5 for the linear correlation coefficient and around 0.2 for the rank correlation coefficient.

4.7 Convergence

It is clearly visible from the out-of-bag curves (<https://github.com/PhilippPro/tuneNtree/tree/master/graphics>) that increasing the number of trees yields a substantial performance gain in most of the cases, but the biggest performance gain in the out-of-bag curves can be seen while growing the first 250 trees. Setting the number of trees from 10 to 250 in the binary classification case provides an average decrease of 0.0306 of the error rate and an increase of 0.0521 of the AUC. On the other hand, using 2000 trees instead of 250 does not yield a big performance gain, the average error rate improvement is only 0.0018 (AUC: 0.0032). The improvement in the multiclass case is bigger with an average improvement of the error rate of 0.0739 (AUC: 0.0665) from 10 trees to 250 and an average improvement of 0.0039 (AUC: 0.0057) for using 2000 trees instead of 250. For regression we have an improvement of 0.1210 of the R^2 within the first 250 trees and an improvement of 0.0039 for using 2000 trees instead of 250. These results are concordant with a comment by Breiman (1996a) (Section 6.2) who notes that fewer bootstrap replicates are necessary when the outcome is numerical and more are required for an increasing number of classes.

5. Conclusions and Extensions

In this section we draw conclusions of the given results and discuss possible extensions.

5.1 Assessment of the Convergence

For the assessment of the convergence in the classification case we generally recommend using measures other than the error rate, such as AUC, the Brier score or the logarithmic loss for which the OOB curves are much more similar as we have seen in our correlation analysis. Their convergence rate is not so dependent on observations with ε_i close to 0.5 (in the binary classification case), and they give an indication of the general stability of the probability estimations of all observations. This can be especially important if the threshold

for classification is not set a priori to 0.5. The new `OOBCurve` R package is a tool to examine the rate of convergence of the trained RF with any measure that is available in the `mlr` R package. It is important to remember that for the calculation of the OOB error curve at T only $\exp(-1) \cdot T$ trees are used. Thus, as far as future independent data is concerned, the convergence of the performances is by $\exp(1) \approx 2.7$ faster than observed from our OOB curves. Having this in mind, our observations (see Section 4.7) are in agreement with the results of Oshiro et al. (2012), who conclude that after growing 128 trees no big gain in the AUC performance could be achieved by growing more trees.

5.2 Why More Trees Are Better

Non-monotonous expected error rate curves observed in the case of binary classification might be seen as an argument in favour of tuning the number T of trees. Our results, however, suggest that tuning is not recommendable in the case of classification. Firstly, non-monotonous patterns are observed only with some performance measures such as the error rate and the AUC in case of classification. Measures such as the Brier score or the logarithmic loss, which are based on probabilities rather than on the predicted class and can thus be seen as more refined, do not yield non-monotonous patterns, as theoretically proved in Section 3 and empirically observed based on a very large number of datasets in Section 4. Secondly, non-monotonous patterns in the expected error rate curves are the result of a particular rare combination of ε_i 's in the training data. Especially if the training dataset is small, the chance is high that the distribution of the ε_i will be different for independent test data, for example values of ε_i close to but larger than 0.5 may not be present. In this case, the expected error rate curve for this independent future dataset would not be non-monotonous, and a large T is better. Thirdly, even in the case of non-monotonic expected error rate curves, the minimal error rate value is usually only slightly smaller than the value at convergence (see Section 4.4.1). We argue that this very small gain - which, as outlined above, is relevant only for future observations with $\varepsilon_i > 0.5$ - probably does not compensate the advantage of using more trees in terms of other performance measures or in terms of the precision of the variable importance measures, which are very commonly used in practice.

In the case of regression, our theoretical results show that the expected out-of-bag mse curve is monotonously decreasing. For the mean absolute error the empirical results suggest the same. In terms of the less common measures *median* squared error and *median* absolute error (as opposed to *mean* losses), however, performance may get worse with increasing number of trees. More research is needed.

5.3 Extensions

Note that our theoretical results are not only valid for random forest but generalizable to any ensemble method that uses a randomization technique, since the fact that the base learners are trees and the specific randomization procedure (for example bagging) do not play any role in our proofs. Our theoretical results could possibly be extended to the multiclass case, as supported by our results obtained with 44 multiclass datasets.

Although we claim that increasing the number of trees cannot harm noticeably as far as measures based on average loss are considered, our empirical results show that for most of the examined datasets, the biggest performance gain is achieved when training the first

100 trees. However, the rate of convergence may be influenced by other hyperparameters of the RF. For example lower sample size while taking bootstrap samples for each tree, bigger constraints on the tree depth or more variables lead to less correlated trees and hence more trees are needed to reach convergence.

One could also think of an automatic break criterion which stops the training automatically according to the convergence of the OOB curves. For example, training could be stopped if the last T_{last} trees did not improve performance by more than Δ , where T_{last} and Δ are parameters that should be fixed by the user as a compromise between performance and computation time. Note that, if variable importances are computed, it may be recommended to also consider their convergence. This issue also requires more research.

Acknowledgments

We would like to thank Alexander Dürre for useful comments on the approximation of the logarithmic loss and Jenny Lee for language editing.

References

- Ranjan Kumar Barman, Sudipto Saha, and Santasabuj Das. Prediction of interactions between viral and host proteins using supervised machine learning methods. *PLOS ONE*, 9(11):1–10, 2014.
- Bernd Bischl, Michel Lang, Lars Kotthoff, Julia Schiffner, Jakob Richter, Erich Studerus, Giuseppe Casalicchio, and Zachary M. Jones. mlr: Machine learning in R. *Journal of Machine Learning Research*, 17(170):1–5, 2016. R package version 2.9.
- Anne-Laure Boulesteix, Rory Wilson, and Alexander Hapfelmeier. Towards evidence-based computational statistics: lessons from clinical research on the role and design of real-data benchmark studies. *BMC Medical Research Methodology*, 17(1):138, 2017.
- Leo Breiman. Bagging predictors. *Machine Learning*, 24(2):123–140, 1996a.
- Leo Breiman. Out-of-bag estimation. *Technical report, Statistics Department, University of California 1996*, 1996b.
- Leo Breiman. Random forests. *Machine Learning*, 45(1):5–32, 2001.
- Giuseppe Casalicchio, Jakob Bossek, Michel Lang, Dominik Kirchoff, Pascal Kerschke, Benjamin Hofner, Heidi Seibold, Joaquin Vanschoren, and Bernd Bischl. OpenML: An R package to connect to the machine learning platform OpenML. *Computational Statistics*, 32(3):1–15, 2017.
- César Ferri, José Hernández-Orallo, and R Modroiu. An experimental comparison of performance measures for classification. *Pattern Recognition Letters*, 30(1):27–38, 2009.
- Jerome H Friedman. On bias, variance, 0/1-loss, and the curse-of-dimensionality. *Data Mining and Knowledge Discovery*, 1(1):55–77, 1997.

- Jerome H. Friedman. Greedy function approximation: a gradient boosting machine. *Annals of Statistics*, 29(5):1189–1232, 2001.
- David J Hand and Robert J Till. A simple generalisation of the area under the ROC curve for multiple class classification problems. *Machine Learning*, 45(2):171–186, 2001.
- Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning*. Springer Series in Statistics. Springer New York Inc., New York, NY, USA, 2001.
- Daniel Hernández-Lobato, Gonzalo Martínez-Muñoz, and Alberto Suárez. How large should ensembles of classifiers be? *Pattern Recognition*, 46(5):1323–1336, 2013.
- Torsten Hothorn, Kurt Hornik, and Achim Zeileis. Unbiased recursive partitioning: A conditional inference framework. *Journal of Computational and Graphical Statistics*, 15(3):651–674, 2006.
- Michel Lang, Bernd Bischl, and Dirk Surmann. batchtools: Tools for R to work on batch systems. *The Journal of Open Source Software*, 2(10), 2017.
- Patrice Latinne, Olivier Debeir, and Christine Decaestecker. Limiting the number of trees in random forests. In *International Workshop on Multiple Classifier Systems*, pages 178–187. Springer, 2001.
- Andy Liaw and Matthew Wiener. Classification and regression by randomForest. *R News*, 2(3):18–22, 2002. R package version 4.6-12.
- Thais Mayumi Oshiro, Pedro Santoro Perez, and José Augusto Baranauskas. How many trees in a random forest? In *International Workshop on Machine Learning and Data Mining in Pattern Recognition*, pages 154–168. Springer, 2012.
- Piergiorgio Palla and Giuliano Armano. *RFmarkerDetector: Multivariate Analysis of Metabolomics Data using Random Forests*, 2016. R package version 1.0.1.
- Arvind Raghu, Praveen Devarsetty, Peiris David, Tarassenko Lionel, and Clifford Gari. Implications of cardiovascular disease risk assessment using the who/ish risk prediction charts in rural india. *PLOS ONE*, 10(8):1–13, 2015.
- Peter J. Rousseeuw. Least median of squares regression. *Journal of the American Statistical Association*, 79(388):871–880, 1984.
- Mark R Segal. Machine learning benchmarks and random forest regression. *Center for Bioinformatics & Molecular Biostatistics*, 2004.
- Carolin Strobl, Anne-Laure Boulesteix, Achim Zeileis, and Torsten Hothorn. Bias in random forest variable importance measures: Illustrations, sources and a solution. *BMC Bioinformatics*, 8(1):25, 2007.
- Joaquin Vanschoren, Jan N. van Rijn, Bernd Bischl, and Luis Torgo. OpenML: Networked science in machine learning. *SIGKDD Explorations*, 15(2):49–60, 2013.
- Marvin N. Wright. *ranger: A Fast Implementation of Random Forests*, 2016. R package version 0.6.0.